

results from DICTION.

Optionally, and preferably, rectangles 83 collectively representing predictions of a single method and/or approach are identically colored and/or textured, and are distinguishable from the color and/or texture used for a different method and/or approach.

Alternatively, or in addition, the color, hue, density, or texture of rectangles 83 can be used further to report a measure of the bioinformatic reliability of the prediction. For example, many gene prediction programs will report a measure of the reliability of prediction. Thus, increasing degrees of such reliability can be indicated, e.g., by increasing density of shading. Where display 80 is used as a graphical user interface, such measures of reliability, and indeed all other results output by the program, can additionally or alternatively be made accessible through linkage from individual rectangles 83, as by time-delayed window ("tool tip" window), or by pointer (e.g., mouse)-activated link.

As earlier described, increased predictive reliability can be achieved by requiring consensus among methods and/or approaches to determining function. Thus, field 81 can include a horizontal series of rectangles 83 that indicate one or more degrees of consensus in predictions of function.

Although FIG. 3 shows three series of horizontally disposed rectangles in field 81, display 80 can include as few as one such series of rectangles and as many as can discriminably be displayed, depending upon the number of methods and/or approaches used to predict a given function.

Furthermore, field 81 can be used to show predictions of a plurality of different functions. However, the increased visual complexity occasioned by such display makes more useful the ability of the user to select

a single function for display. When display 80 is used as a graphical user interface for computer query and analysis, such function can usefully be indicated and user-selectable, as by a series of graphical buttons or tabs
5 (not shown in FIG. 3).

Rectangle 89 is shown in FIG. 3 as including interposed rectangle 84. Rectangle 84 represents the portion of annotated sequence for which predicted functional information has been assayed physically, with
10 the starting and ending nucleotides of the assayed material indicated by the X axis coordinates of the left and right borders of rectangle 84. Rectangle 85, with optional inclusive circles 86 (86a, 86b, and 86c) displays the results of such physical assay.

15 Although a single rectangle 84 is shown in FIG. 3, physical assay is not limited to just one region of annotated genomic sequence. It is expected that an increasing percentage of regions predicted to have function by process 200 will be assayed physically, and that display
20 80 will accordingly, for any given genomic sequence, have an increasing number of rectangles 84 and 85, representing an increased density of sequence annotation.

Where the function desired to be identified is protein coding, rectangle 84 identifies the sequence of the
25 probe used to measure expression. In embodiments of the present invention where expression is measured using genome-derived single exon microarrays, rectangle 84 identifies the sequence included within the probe immobilized on the support surface of the microarray. As
30 noted *supra*, such probe will often include a small amount of additional, synthetic, material incorporated during amplification and designed to permit reamplification of the probe, which sequence is typically not shown in display 80.

Rectangle 87 is used to present the results of
35 bioinformatic assay of the genomic sequence. For example,

where the function desired to be identified is protein coding, process 400 can include bioinformatic query of expression databases with the sequences predicted in process 200 to encode exons. And as earlier discussed, because bioinformatic assay presents fewer constraints than does physical assay, often the entire output of process 200 can be used for such assay, without further subsetting thereof by process 300. Therefore, rectangle 87 typically need not have separate indicators therein of regions submitted for bioinformatic assay; that is, rectangle 87 typically need not have regions therein analogous to rectangles 84 within rectangle 89.

Rectangle 87 as shown in FIG. 3 includes smaller rectangles 880 and 88. Rectangles 880 indicate regions that returned a positive result in the bioinformatic assay, with rectangles 88 representing regions that did not return such positive results. Where the function desired to be predicted and displayed is protein coding, rectangles 880 indicate regions of the predicted exons that identify sequence with significant similarity in expression databases, such as EST, SNP, SAGE databases, with rectangles 88 indicating genes novel over those identified in existing expression data bases.

Rectangles 880 can further indicate, through color, shading, texture, or the like, additional information obtained from bioinformatic assay.

For example, where the function assayed and displayed is protein coding, the degree of shading of rectangles 880 can be used to represent the degree of sequence similarity found upon query of expression databases. The number of levels of discrimination can be as few as two (identity, and similarity, where similarity has a user-selectable lower threshold). Alternatively, as many different levels of discrimination can be indicated as can visually be discriminated.

Where display 80 is used as a graphical user interface, rectangles 880 can additionally provide links directly to the sequences identified by the query of expression databases, and/or statistical summaries thereof.

5 As with each of the precedingly-discussed uses of display 80 as a graphical user interface, it should be understood that the information accessed via display 80 need not be resident on the computer presenting such display, which often will be serving as a client, with the linked
10 information resident on one or more remotely located servers.

Rectangle 85 displays the results of physical assay of the sequence delimited by its left and right borders.

15 Rectangle 85 can consist of a single rectangle, thus indicating a single assay, or alternatively, and increasingly typically, will consist of a series of rectangles (85a, 85b, 85c) indicating separate physical assays of the same sequence.

20 Where the function assayed is gene expression, and where gene expression is assayed as herein described using simultaneous two-color fluorescent detection of hybridization to genome-derived single exon microarrays, individual rectangles 85 can be colored to indicate the
25 degree of expression relative to control. Conveniently, shades of green can be used to depict expression in the sample over control values, and shades of red used to depict expression less than control, corresponding to the spectra of the Cy3 and Cy5 dyes conventionally used for
30 respective labeling thereof. Additional functional information can be provided in the form of circles 86 (86a, 86b, 86c), where the diameter of the circle can be used to indicate expression intensity. As discussed *infra*, such relative expression (expression ratios) and absolute
35 expression (signal intensity) can be expressed using

normalized values.

Where display 80 is used as a graphical user interface, rectangle 85 can be used as a link to further information about the assay. For example, where the assay
5 is one for gene expression, each rectangle 85 can be used to link to information about the source of the hybridized mRNA, the identity of the control, raw or processed data from the microarray scan, or the like.

FIG. 4 is rendition of display 80 representing
10 gene prediction and gene expression for a hypothetical BAC, showing conventions used in the Examples presented *infra*. BAC sequence ("Chip seq.") 89 is presented, with the physically assayed region thereof (corresponding to rectangle 84 in FIG. 3) shown in white. Algorithmic gene
15 predictions are shown in field 81, with predictions by GRAIL shown, predictions by GENEFINDER, and predictions by DICTION shown. Within rectangle 87, regions of sequence that, when used to query expression databases, return identical or similar sequences ("EST hit") are shown as
20 white rectangles (corresponding to rectangles 880 in FIG. 3), gray indicates low homology, and black indicates unknowns (where black and gray would correspond to rectangles 88 in FIG. 3).

Although FIGS. 3 and 4 show a single stretch of
25 sequence, uninterrupted from left to right, longer sequences are usefully represented by vertical stacking of such individual Mondrians, as shown in FIGS. 9 and 10.

Single Exon Probes Useful For Measuring Gene Expression

30

The methods and apparatus of the present invention rapidly produce functional information from genomic sequence. Where the function to be identified is protein coding, the methods and apparatus of the present
35 invention rapidly identify and confirm the expression of

portions of genomic sequence that function to encode protein. As a direct result, the methods and apparatus of the present invention rapidly yield large numbers of single-exon nucleic acid probes, the majority from
5 previously unknown genes, each of which is useful for measuring and/or surveying expression of a specific gene in one or more tissues or cell types.

It is, therefore, another aspect of the present invention to provide genome-derived single exon nucleic
10 acid probes useful for gene expression analysis, and particularly for gene expression analysis by microarray.

Using the methods and genome-derived single-exon microarrays of the present invention, we have for example readily identified a large number of unique ORFs from human
15 genomic sequence. Using single exon probes that encompass these ORFs, we have demonstrated, through microarray hybridization analysis, the expression of 9,980 of these ORFs in heart.

As would immediately be appreciated by one of
20 skill in the art, each single exon probe having demonstrable expression in heart is currently available for use in measuring the level of its ORF's expression in heart.

Diseases of the heart and vascular system are a
25 significant cause of human morbidity and mortality. Increasingly, genetic factors are being found that contribute to predisposition, onset, and/or aggressiveness of most, if not all, of these diseases. Although mutations in single genes have on occasion been identified as
30 causative, these disorders are for the most part believed to have polygenic etiologies.

For example, cardiovascular disease (CVD), which includes coronary heart disease, stroke, and peripheral arterial vascular disease, is the leading cause of death in
35 the United States and other developed countries. In

developing regions, coronary heart disease and stroke are ranked second and third, respectively, as causes of mortality. In the United States alone, about 1 million deaths (about 42% of total deaths per year) result from CVD
5 each year. CVD is also a significant cause of morbidity, with about 1.5 million people suffering myocardial infarction, and about 500,000 suffering strokes in the United States each year. With risk for CVD increasing with age, and an increasingly aging population, CVD will
10 continue to be a major health problem into the future.

CVD is caused by arterial lesions that begin as fatty streaks, which consist of lipid-laden foam cells, and develop into fibrous plaques. The atherosclerotic plaque may grow slowly, and over several decades may produce a
15 severe stenosis or result in arterial occlusion. Some plaques are stable, but other, more unstable, ones may rupture and induce thrombosis. The thrombi may embolize, rapidly occluding the lumen and leading to myocardial infarction or acute ischemic syndrome.

Risk factors for CVD include age and gender. In addition, a family history of CVD significantly increases risk, indicating a genetic basis for development of this disease complex. Obesity, especially truncal obesity, the cause of which is suspected to be genetic, is yet another
25 risk factor for CVD. Familial disorders such as hyperlipidemia, hypoalphalipoproteinemia, hypertriglyceridemia, hypercholesterolemia, hyperinsulinemia, homocystinuria, and dysbetalipoproteinemia, all of which lead to lipid or
30 lipoprotein abnormalities, can predispose one to the development of CVD. Both insulin-dependent and non-insulin-dependent diabetes mellitus, both of which have genetic components, have been also linked to the development of atherosclerosis.

35 The literature is replete with evidence for

genetic causes of cardiovascular diseases. For example, studies by Allayee et al., Am. J. Hum. Genet. 63:577-585(1998), indicated a genetic association between familial combined hyperlipidemia (FCHL) and small dense LDL particles. The studies also concluded that the genetic determinants for LDL particle size are shared, at least in part, among FCHL families and the more general population at risk for CVD. Juo et al., Am. J. Hum. Genet. 63: 586-594 (1998) demonstrated that small, dense LDL particles and elevated apolipoprotein B levels, both of which are commonly found in members of FCHL families, share a common major gene plus individual polygenic components.

The common major gene was estimated to explain 37% of the variants of adjusted LDL particle size and 23% of the variants of adjusted apoB levels.

The atherogenic lipoprotein phenotype (ALP) is a common heritable trait, symptoms of which include a prevalence of small, dense LDL particles, increased levels of triglyceride-rich lipoproteins, reduced levels of high density lipoprotein, and increased risk of CVD, particularly myocardial infarction. Both Nishina et al., Proc. Nat. Acad. Sci. 89: 708-712 (1992) and Rotter et al., Am. J. Hum. Genet. 58: 585-594(1996) demonstrated linkage between ALP and the LDLR locus. Rotter et al., supra, also reported linkage to the CETP locus on chromosome 16 and to the SOD1 locus on chromosome 6, and possibly also to the APOA1/APOC3/APOA4 cluster on chromosome 11.

Mutations in genes identified as components of lipid metabolism, e.g., apolipoprotein E (apoE) and LDL receptor (LDLR), have been shown to be associated with predisposition to the development of CVD. For example, several apoE variants had been found to be associated with familial dysbetalipoproteinemia, characterized by elevated plasma cholesterol and triglyceride levels and an increased risk for atherosclerosis (de Knijff et al., Mutat 4: 178-

194 (1994)). Mutations in the LDLR gene have been associated with the familial hypercholesterolemia, an autosomal dominant disorder characterized by elevation of serum cholesterol bound to low density lipoprotein (LDL),
5 that can lead to increased susceptibility to CVD.

To date, mutations in numerous genes have been shown to be associated with increased CVD susceptibility. However, the identified genetic associations are believed not to account for all genetic contributions to CVD.

10 As yet another example, hypertension is a major health problem because of its high prevalence and its association with increased risk of CVD. Approximately 25% of all adults and over 60% of persons older than 60 years in the United States have high blood pressure.

15 Arterial or systemic hypertension is diagnosed when the average of two or more diastolic BP measurements on at least two subsequent visits is 90 mm Hg or more, or when the average of multiple systolic BP readings on two or more subsequent visits is consistently greater than 140 mm
20 Hg. Pulmonary hypertension is defined as pressure within the pulmonary arterial system elevated above the normal range; pulmonary hypertension may lead to right ventricle (RV) failure.

Hypertension, together with other cardiovascular
25 risk factors, leads to atherosclerosis and other forms of CVD, primarily by damaging the vascular endothelium. In more than 40% of the U.S. population, hypertension is accompanied by hyperlipidemia and leads to the development of atherosclerotic plaques. In the absence of
30 hyperlipidemia, intimal thickening occurs. Non-atherosclerotic hypertension-induced vascular damage can lead to stroke or heart failure.

Familial diseases associated with secondary hypertension include familial renal disease, polycystic
35 kidney disease, medullary thyroid cancer, pheochromocytoma,

and hyperparathyroidism. Hypertension is also twice as common in patients with diabetes mellitus.

More than 95% of all hypertension cases are essential hypertension, that is, lack identifiable
5 antecedent clinical cause. Essential hypertension shows clustering in families and can result from a variety of genetic diseases. In most cases, high blood pressure results from a complex interaction of factors with both genetic and environmental components. The recent search
10 for genes that contribute to the development of essential hypertension has shown that the disorder is polygenic in origin. However, with several exceptions (such as angiotensinogen, angiotensin receptor-1, beta-3 subunit of guanine nucleotide-binding protein, tumor necrosis factor
15 receptor-2, and "-adducin), the particular genes involved are still being sought.

Susceptibility loci for essential hypertension have been mapped to chromosomes 17 and 15q. Hasstedt et al., Am. J. Hum. Genet. 43: 14-22 (1988) measured red cell
20 sodium in 1,800 normotensive members of 16 Utah pedigrees ascertained through hypertensive or normotensive probands, siblings with early stroke death, or brothers with early coronary disease, and suggested that red blood cell sodium was determined by 4 alleles at a single locus. This major
25 locus was thought to explain 29% of the variance in red cell sodium, and polygenic inheritance explained another 54.6%. A higher frequency of the high red blood cell sodium genotype in pedigrees in which the proband was hypertensive rather than normotensive provided evidence that this major
30 locus increases susceptibility to hypertension.

From a study of systolic blood pressure in 278 pedigrees, Perusse et al., Am. J. Hum. Genet. 49: 94-105 (1991) reported that variability in systolic blood pressure is likely influenced by allelic variation of a single gene,
35 with gender and age dependence. They also suggested that a

single gene may be associated with a steeper increase of blood pressure with age among males and females.

There is strong evidence, however, for additional as yet uncharacterized, hypertension-associated loci on
5 other chromosomes.

For example, Xu et al., Am. J. Hum. Genet. 64: 1694-1701 (1999) carried out a systematic search for chromosomal regions containing genes that regulate blood pressure by scanning the entire autosomal genome using 367
10 polymorphic markers. Because of the sampling design, the number of sib pairs, and the availability of genotyped parents, this study represented one of the most powerful of its kind. Although no regions achieved a 5% genomewide significance level, maximum lod scores were greater than
15 2.0 for regions of chromosomes 3, 11, 15, 16, and 17.

As another example, cardiac arrhythmias account for several thousand deaths each year. Arrhythmias such as ventricular fibrillation, which causes more than 300,000 sudden deaths annually in the United States alone,
20 encompass a multitude of disorders. Another type of arrhythmia, idiopathic dilated cardiomyopathy, of which familial dilated cardiomyopathy accounts for 20-25%, is responsible for more than 10,000 deaths in the United States annually and is the predominant indication for
25 cardiac transplantation.

Cardiac arrhythmias can be divided into bradyarrhythmias (slowed rhythms) or tachyarrhythmias (speeded rhythms). Bradyarrhythmias result from abnormalities of intrinsic automatic behavior or
30 conduction, primarily within the atrioventricular node and the His-Purkinje's network. Tachyarrhythmias are caused by altered automaticity, reentry, or triggered automaticity.

Bradyarrhythmias arising from suspected polygenic disorders include Long QT syndrome 4, atrioventricular
35 block, familial sinus node disease, progressive cardiac

conduction defect, and familial cardiomyopathy.

Tachyarrhythmias with possible underlying polygenic causes include familial ventricular tachycardia, Wolff-Parkinson-White syndrome, familial arrhythmogenic right ventricular dysplasia, heart-hand syndrome V, Mal de Meleda, familial
5 ventricular fibrillation, and familial noncompaction of left ventricular myocardium.

For some of the arrhythmias, one or more of the causative genes have been identified.

10 For example, atrioventricular block has been associated with mutations in the SCN5A gene, as well as mutations in a locus mapped to 19q13. Studies have shown linkage of familial sinus node disease to a marker on 10q22-q24. Familial ventricular tachycardia has been
15 linked to mutations in genes encoding the G protein subunit alpha-i2 (GNAI1), and/or related genes. Examination of families with Wolff-Parkinson-White syndrome suggest an autosomal dominant pattern of inheritance and evidence of linkage of the disorder to DNA markers on band 7q3.
20 Linkage analysis shows strong evidence for localization of a gene for Mal de Meleda disease on 8qter. Familial ventricular fibrillation can be caused by mutations in the cardiac sodium channel gene SCN5A. Familial noncompaction of left ventricular myocardium has been linked to mutations
25 in the gene encoding tafazzin (TAZ), or in the FK506-binding protein 1A gene (FKBP1A).

Familial dilated cardiomyopathy is characterized by an autosomal dominant pattern of inheritance with age-related penetrance. The linkage of familial dilated
30 cardiomyopathy to several loci indicate that it is polygenic. These loci include CMD1A on 1p11-q11, CMD1B on 9q13, CMD1C on 10q21, CMD1D on 1q32, CMD1E on 3p, CMD1F on 6q, CMD1G on 2q31, CMD1H on 2q14-q22, and CMD1I, which results from mutation in the DES gene on 2q35.

35 In addition, cardiomyopathy can also be caused by

mutations in the ACTC gene, the cardiac beta-myosin heavy chain gene (MYH7), or the cardiac troponin T gene.

Familial arrhythmogenic right ventricular dysplasia is inherited as an autosomal dominant with
5 reduced penetrance and is one of the major genetic causes of juvenile sudden death. It is estimated that the prevalence of familial arrhythmogenic right ventricular dysplasia ranges from 6 per 10,000 in the general population to 4.4 per 1,000 in some areas.

10 Several loci for familial arrhythmogenic right ventricular dysplasia have been mapped indicating that this disease is also polygenic in nature. These loci include ARVD1 on 14q23-q24, ARVD2 on 1q42-q43, ARVD3 on 14q12-q22, ARVD4 on 2q32.1-q32.3, ARVD5 on 3p23, and ARVD6 on 10p14-
15 p12.

Progressive cardiac conduction defect (PCCD), also called Lenegre-Lev disease, is one of the most common cardiac conduction diseases. It is characterized by progressive alteration of cardiac conduction through the
20 His-Purkinje system with right or left bundle branch block and widening of QRS complexes, leading to complete atrioventricular block and ultimately causing syncope and sudden death. It represents the major cause of pacemaker implantation in the world (0.15 implantations per 1,000
25 inhabitants per year in developed countries). The cause of PCCD is unknown but familial cases with right bundle branch block have been reported suggesting that at least some cases are of genetic origin. Reports have linked PCCD to HB1 on 19q13.3, and to mutations in the SCN5A gene (Schott
30 et al., Nature Genet. 23: 20-21 (1999)).

As yet a further example, congenital heart disease occurs at a rate of 8 per 1000 live births, which corresponds to approximately 32,000 infants with newly diagnosed congenital heart disease each year in the United
35 States. Twenty percent of infants with congenital heart

disease die within the first year of life. Approximately 80% of the first-year survivors live to reach adulthood. Congenital heart disease also has economic impact due to the estimated 20,000 surgical procedures performed to
5 correct circulatory defects in these patients. The estimated number of adults with congenital heart disease in the United States is currently about 900,000.

In 90% of patients, congenital heart disease is attributable to multifactorial inheritance. Only 5-10% of
10 malformations are due to primary genetic factors, which are either chromosomal or a result of a single mutant gene.

The most common congenital heart disease found in adults is bicuspid aortic valve. This defect occurs in 2% of the general population and accounts for approximately
15 50% of operated cases of aortic stenosis in adults. Atrial septal defect is responsible for 30-40% of congenital heart disease seen in adults. The most common congenital cardiac defect observed in the pediatric population is ventricular septal defect, which accounts for 15-20% of all congenital
20 lesions. Tetralogy of Fallot is the most common cyanotic congenital anomaly observed in adults. Other congenital heart diseases include Eisenmenger's syndrome, patent ductus arteriosus, pulmonary stenosis, coarctation of the aorta, transposition of the great arteries, tricuspid
25 atresia, univentricular heart, Ebstein's anomaly, and double-outlet right ventricle.

A number of studies have identified putative genetic loci associated with one or more congenital heart diseases.

30 Congenital heart disease affects more than 40% of all Down syndrome patients. The candidate chromosomal region containing the putative gene or genes for congenital heart disease associated with Down syndrome is 21q22.2-q22.3, between ETS2 and MX1.

35 DiGeorge syndrome (DGS) is characterized by

several symptoms including outflow tract defects of the heart such as teratology of Fallot. Most cases result from a deletion of chromosome 22q11.2 (the DiGeorge syndrome chromosome region, or DGCR). The 22q11 deletion is the
5 second most common cause of congenital heart disease after Down syndrome. Several genes are lost in this deletion including the putative transcription factor TUPLE1. This deletion is associated with a variety of phenotypes, e.g., Shprintzen syndrome; conotruncal anomaly face (or Takao
10 syndrome); and isolated outflow tract defects of the heart including Tetralogy of Fallot, truncus arteriosus, and interrupted aortic arch.

Whereas 90% of cases of DGS may now be attributed to a 22q11 deletion, other associated chromosome defects
15 have been identified. For example, Greenberg et al., Am. J. Hum. Genet. 43:605-611 (1988), reported 1 case of DGS with del10p13 and one with a 18q21.33 deletion. Fukushima et al., Am. J. Hum. Genet. 51 (suppl.):A80 (1992) reported linkage with a deletion of 4q21.3-q25. Gottlieb et al.,
20 Am. J. Hum. Genet. 62: 495-498 (1998) concluded that the deletion of more than 1 region on 10p could be associated with the DGS phenotype. The association of the DiGeorge syndrome with at least 2 and possibly more chromosomal locations suggests strongly the involvement of several
25 genes in this disease.

Digilio et al., J. Med. Genet. 34: 188-190 (1997), calculated empiric risk figures for recurrence of isolated Tetralogy of Fallot in families after exclusion of del(22q11), and concluded that gene(s) different from those
30 located on 22q11 must be involved in causing familial aggregation of nonsyndromic Tetralogy of Fallot. Johnson et al., Am. J. Med. Genet. (1997) conducted a cytogenetic evaluation of 159 cases of Tetralogy of Fallot. They reported that a del(22q11) was identified in 14% who
35 underwent fluorescence in situ hybridization (FISH) testing

with the N25 cosmid probe.

Other congenital heart disease are also suspected to be of polygenic origin. For example, Holmes et al., Birth Defects Orig. Art. Ser. X(4): 228-230 (1974) described familial clustering of hypoplastic left heart syndrome in siblings consistent with multifactorial causation.

Other significant diseases of the heart and vascular system are also believed to have a genetic, typically polygenic, etiological component. These diseases include, for example, hypoplastic left heart syndrome, cardiac valvular dysplasia, Pfeiffer cardiocranial syndrome, oculofaciocardiodental syndrome, Kapur-Toriello syndrome, Sonoda syndrome, Ohdo Blepharophimosis syndrome, heart-hand syndrome, Pierre-Robin syndrome, Hirschsprung disease, Kousseff syndrome, Grange occlusive arterial syndrome, Kearns-Sayre syndrome, Kartagener syndrome, Alagille syndrome, Ritscher-Schinzel syndrome, Ivemark syndrome, Young-Simpson syndrome, hemochromatosis, Holzgreve syndrome, Barth syndrome, Smith-Lemli-Opitz syndrome, glycogen storage disease, Gaucher-like disease, Fabry disease, Lowry-Maclean syndrome, Rett syndrome, Opitz syndrome, Marfan syndrome, Miller-Dieker lissencephaly syndrome, mucopolysaccharidosis, Bruada syndrome, humerospinal dysostosis, Phaver syndrome, McDonough syndrome, Marfanoid hypermobility syndrome, atransferrinemia, Cornelia de Lange syndrome, Leopard syndrome, Diamond-Blackfan anemia, Steinfeld syndrome, progeria, and Williams-Beuren syndrome.

The human genome-derived single exon nucleic acid probes and microarrays of the present invention are useful for predicting, diagnosing, grading, staging, monitoring and prognosing diseases of human heart and vascular system, particularly those diseases with polygenic etiology. With each of the single exon probes described herein shown to be

expressed at detectable levels in human heart, and with about 2/3 of the probes identifying novel genes, the single exon microarrays of the present invention provide exceptionally high informational content for such studies.

5 For example, diagnosis (including differential diagnosis among clinically indistinguishable disorders), staging, and/or grading of a disease can be based upon the quantitative relatedness of a patient gene expression profile to one or more reference expression profiles known
10 to be characteristic of a given heart or vascular disease, or to specific grades or stages thereof.

 In one embodiment, the patient gene expression profile is generated by hybridizing nucleic acids obtained directly or indirectly from transcripts expressed in the
15 patient's heart or vascular tissues to the genome-derived single exon microarray of the present invention. Reference profiles are obtained similarly by hybridizing nucleic acids obtained directly or indirectly from transcripts expressed in heart or vascular tissue of individuals with
20 known disease. Methods for quantitatively relating gene expression profiles, without regard to the function of the protein encoded by the gene, are disclosed in WO 99/58720, incorporated herein by reference in its entirety.

 In another approach, the genome-derived single
25 exon probes and microarrays of the present invention can be used to interrogate genomic DNA, rather than pools of expressed message; this latter approach permits predisposition to and/or prognosis of heart or vascular disease to be assessed through the massively parallel
30 determination of altered copy number, deletion, or mutation in the patient's genome of exons known to be expressed in human heart. The algorithms set forth in WO 99/58720 can be applied to such genomic profiles without regard to the function of the protein encoded by the interrogated gene.

35 The utility is specific to the probe; at

sufficiently high hybridization stringency, which stringencies are well known in the art — see Ausubel et al. and Maniatis et al. — each probe reports the level of expression of message specifically containing that ORF.

5 It should be appreciated, however, that the probes of the present invention, for which expression in the heart has been demonstrated are useful for both measurement in the heart and for survey of expression in other tissues.

10 Significant among such advantages is the presence of probes for novel genes.

 As mentioned above and further detailed in Examples 1 and 2, the methods described enable ORFs which are not present in existing expression databases to be
15 identified. And the fewer the number of tissues in which the ORF can be shown to be expressed, the more likely the ORF will prove to be part of a novel gene: as further discussed in Example 2, ORFs whose expression was measurable in only a single of the tested tissues were
20 represented in existing expression databases at a rate of only 11%, whereas 36% of ORFs whose expression was measurable in 9 tissues were present in existing expression databases, and fully 45% of those ORFs expressed in all ten tested tissues were present in existing expressed sequence
25 databases.

 Either as tools for measuring gene expression or tools for surveying gene expression, the genome-derived single exon probes of the present invention have significant advantages over the cDNA or EST-based probes
30 that are currently available for achieving these utilities.

 The genome-derived single exon probes of the present invention are useful in constructing genome-derived single exon microarrays; the genome-derived single exon microarrays, in turn, are useful devices for measuring and
35 for surveying gene expression in the human.

Gene expression analysis using microarrays – conventionally using microarrays having probes derived from expressed message – is well-established as useful in the biological research arts (see Lockhart et al. *Nature* 405, 5 827-836).

Microarrays have been used to determine gene expression profiles in cells in response to drug treatment (see, for example, Kaminski et al., "Global Analysis of Gene Expression in Pulmonary Fibrosis Reveals Distinct 10 Programs Regulating Lung Inflammation and Fibrosis," *Proc. Natl. Acad. Sci. USA* 97(4):1778-83 (2000); Bartosiewicz et al., "Development of a Toxicological Gene Array and Quantitative Assessment of This Technology," *Arch. Biochem. Biophys.* 376(1):66-73 (2000)), viral infection (see for 15 example, Geiss et al., "Large-scale Monitoring of Host Cell Gene Expression During HIV-1 Infection Using cDNA Microarrays," *Virology* 266(1):8-16 (2000)) and during cell processes such as differentiation, senescence and apoptosis (see, for example, Shelton et al., "Microarray Analysis of 20 Replicative Senescence," *Curr. Biol.* 9(17):939-45 (1999); Voehringer et al., "Gene Microarray Identification of Redox and Mitochondrial Elements That Control Resistance or Sensitivity to Apoptosis," *Proc. Natl. Acad. Sci. USA* 97(6):2680-5 (2000)).

Microarrays have also been used to determine abnormal gene expression in diseased tissues (see, for example, Alon et al., "Broad Patterns of Gene Expression Revealed by Clustering Analysis of Tumor and Normal Colon Tissues Probed by Oligonucleotide Arrays," *Proc. Natl. Acad. Sci. USA* 96(12):6745-50 (1999); Perou et al., 30 "Distinctive Gene Expression Patterns in Human Mammary Epithelial Cells and Breast Cancers," *Proc. Natl. Acad. Sci. USA* 96(16):9212-7 (1999); Wang et al., "Identification of Genes Differentially Over-expressed in Lung Squamous Cell 35 Carcinoma Using Combination of cDNA Subtraction and

Microarray Analysis," *Oncogene* 19(12):1519-28 (2000);
Whitney *et al.*, "Analysis of Gene Expression in Multiple
Sclerosis Lesions Using cDNA Microarrays," *Ann. Neurol.*
46(3):425-8 (1999)), in drug discovery screens (see, for
5 example, Scherf *et al.*, "A Gene Expression Database for the
Molecular Pharmacology of Cancer," *Nat. Genet.* 24(3):236-44
(2000)) and in diagnosis to determine appropriate treatment
strategies (see, for example, Sgroi *et al.*, "In vivo Gene
Expression Profile Analysis of Human Breast Cancer
10 Progression," *Cancer Res.* 59(22):5656-61 (1999)).

In microarray-based gene expression screens of
pharmacological drug candidates upon cells, each probe
provides specific useful data. In particular, it should be
appreciated that even those probes that show no change in
15 expression are as informative as those that do change,
serving, in essence, as negative controls.

For example, where gene expression analysis is
used to assess toxicity of chemical agents on cells, the
failure of the agent to change a gene's expression level is
20 evidence that the drug likely does not affect the pathway
of which the gene's expressed protein is a part.
Analogously, where gene expression analysis is used to
assess side effects of pharmacological agents — whether in
lead compound discovery or in subsequent screening of lead
25 compound derivatives — the inability of the agent to alter
a gene's expression level is evidence that the drug does
not affect the pathway of which the gene's expressed
protein is a part.

WO 99/58720 provides methods for quantifying the
30 relatedness of a first and second gene expression profile
and for ordering the relatedness of a plurality of gene
expression profiles. The methods so described permit
useful information to be extracted from a greater
percentage of the individual gene expression measurements
35 from a microarray than methods previously used in the art.

Other uses of microarrays are described in Gerhold *et al.*, *Trends Biochem. Sci.* 24(5):168-173 (1999) and Zweiger, *Trends Biotechnol.* 17(11):429-436 (1999); Schena *et al.*

5 The invention particularly provides genome-derived single-exon probes known to be expressed in heart.

 The individual single exon probes can be provided in the form of substantially isolated and purified nucleic acid, typically, but not necessarily, in a quantity
10 sufficient to perform a hybridization reaction.

 Such nucleic acid can be in any form directly hybridizable to the message that contains the probe's ORF, such as double stranded DNA, single-stranded DNA complementary to the message, single-stranded RNA
15 complementary to the message, or chimeric DNA/RNA molecules so hybridizable. The nucleic acid can alternatively or additionally include either nonnative nucleotides, alternative internucleotide linkages, or both, so long as complementary binding can be obtained. For example, probes
20 can include phosphorothioates, methylphosphonates, morpholino analogs, and peptide nucleic acids (PNA), as are described, for example, in U.S. Patent Nos. 5,142,047; 5,235,033; 5,166,315; 5,217,866; 5,184,444; 5,861,250.

 Usefully, however, such probes are provided in a
25 form and quantity suitable for amplification, where the amplified product is thereafter to be used in the hybridization reactions that probe gene expression. Typically, such probes are provided in a form and quantity suitable for amplification by PCR or by other well known
30 amplification technique. One such technique additional to PCR is rolling circle amplification, as is described, *inter alia*, in U.S. Patent Nos. 5,854,033 and 5,714,320 and international patent publications WO 97/19193 and WO 00/15779. As is well understood, where the probes are
35 to be provided in a form suitable for amplification, the

range of nucleic acid analogues and/or internucleotide linkages will be constrained by the requirements and nature of the amplification enzyme.

Where the probe is to be provided in form
5 suitable for amplification, the quantity need not be sufficient for direct hybridization for gene expression analysis, and need be sufficient only to function as an amplification template, typically at least about 1, 10 or 100 pg or more.

10 Each discrete amplifiable probe can also be packaged with amplification primers, either in a single composition that comprises probe template and primers, or in a kit that comprises such primers separately packaged therefrom. As earlier mentioned, the ORF-specific
15 5' primers used for genomic amplification can have a first common sequence added thereto, and the ORF-specific 3' primers used for genomic amplification can have a second, different, common sequence added thereto, thus permitting, in this embodiment, the use of a single set of 5' and 3'
20 primers to amplify any one of the probes. The probe composition and/or kit can also include buffers, enzyme, etc., required to effect amplification.

As mentioned earlier, when intended for use on a genome-derived single exon microarray of the present
25 invention, the genome-derived single exon probes of the present invention will typically average at least about 100, 200, 300, 400 or 500 bp in length, including (and typically, but not necessarily centered about) the ORF. Furthermore, when intended for use on a genome-derived
30 single exon microarray of the present invention, the genome-derived single exon probes of the present invention will typically not contain a detectable label.

When intended for use in solution phase hybridization, however — that is, for use in a
35 hybridization reaction in which the probe is not first

bound to a support substrate (although the target may indeed be so bound) – length constraints that are imposed in microarray-based hybridization approaches will be relaxed, and such probes will typically be labeled.

5 In such case, the only functional constraint that dictates the minimum size of such probe is that each such probe must be capable of specifically identifying in a hybridization reaction the exon from which it is drawn. In theory, a probe of as little as 17 nucleotides is capable
10 of uniquely identifying its cognate sequence in the human genome. For hybridization to expressed message – a subset of target sequence that is much reduced in complexity as compared to genomic sequence – even fewer nucleotides are required for specificity.

15 Therefore, the probes of the present invention can include as few as 20, 25 or 50 bp or ORF, or more. In particular embodiments, the ORF sequences are given in SEQ ID NOS. 9,981 – 19,771, respectively, for probe SEQ ID NOS. 1 – 9,980. The minimum amount of ORF required to be
20 included in the probe of the present invention in order to provide specific signal in either solution phase or microarray-based hybridizations can readily be determined for each of ORF SEQ ID NOS. 9,981 – 19,771 individually by routine experimentation using standard high stringency
25 conditions.

 Such high stringency conditions are described, *inter alia*, in Ausubel et al. and Maniatis et al. For microarray-based hybridization, standard high stringency conditions can usefully be 50% formamide, 5X SSC, 0.2 µg/µl
30 poly(dA), 0.2 µg/µl human c₀t1 DNA, and 0.5 % SDS, in a humid oven at 42°C overnight, followed by successive washes of the microarray in 1X SSC, 0.2% SDS at 55°C for 5 minutes, and then 0.1X SSC, 0.2% SDS, at 55°C for 20 minutes. For solution phase hybridization, standard high
35 stringency conditions can usefully be aqueous hybridization

at 65°C in 6X SSC. Lower stringency conditions, suitable for cross-hybridization to mRNA encoding structurally- and functionally-related proteins, can usefully be the same as the high stringency conditions but with reduction in
5 temperature for hybridization and washing to room temperature (approximately 25°C).

When intended for use in solution phase hybridization, the maximum size of the single exon probes of the present invention is dictated by the proximity of
10 other expressed exons in genomic DNA: although each single exon probe can include intergenic and/or intronic material contiguous to the ORF in the human genome, each probe of the present invention will include portions of only one expressed exon.

15 Thus, each single exon probe will include no more than about 25 kb of contiguous genomic sequence, more typically no more than about 20 kb of contiguous genomic sequence, more usually no more than about 15 kb, even more usually no more than about 10 kb. Usually, probes that are
20 maximally about 5 kb will be used, more typically no more than about 3 kb.

It will be appreciated that the Sequence Listing appended hereto presents, by convention, only that strand of the probe and ORF sequence that can be directly
25 translated reading from 5' to 3' end. As would be well understood by one of skill in the art, single stranded probes must be complementary in sequence to the ORF as present in an mRNA; it is well within the skill in the art to determine such complementary sequence. It will further
30 be understood that double stranded probes can be used in both solution-phase hybridization and microarray-based hybridization if suitably denatured.

Thus, it is an aspect of the present invention to provide single-stranded nucleic acid probes that have
35 sequence complementary to those described herein above and

below, and double-stranded probes one strand of which has sequence complementary to the probes described herein.

The probes can, but need not, contain intergenic and/or intronic material that flanks the ORF, on one or both sides, in the same linear relationship to the ORF that the intergenic and/or intronic material bears to the ORF in genomic DNA. The probes do not, however, contain nucleic acid derived from more than one expressed ORF.

And when intended for use in solution hybridization, the probes of the present invention can usefully have detectable labels. Nucleic acid labels are well known in the art, and include, *inter alia*, radioactive labels, such as ^3H , ^{32}P , ^{33}P , ^{35}S , ^{125}I , ^{131}I ; fluorescent labels, such as Cy3, Cy5, Cy5.5, Cy7, SYBR[®]

Green and other labels described in Haugland, *Handbook of Fluorescent Probes and Research Chemicals*, 7th ed., Molecular Probes Inc., Eugene, OR (2000), or fluorescence resonance energy transfer tandem conjugates thereof; labels suitable for chemiluminescent and/or enhanced chemiluminescent detection; labels suitable for ESR and NMR detection; and labels that include one member of a specific binding pair, such as biotin, digoxigenin, or the like.

The probes, either in quantity sufficient for hybridization or sufficient for amplification, can be provided in individual vials or containers.

Alternatively, such probes can usefully be packaged as a plurality of such individual genome-derived single exon probes.

When provided as a collection of plural individual probes, the probes are typically made available in amplifiable form in a spatially-addressable ordered set, typically one per well of a microtiter dish. Although a 96 well microtiter plate can be used, greater efficiency is obtained using higher density arrays.

If, as earlier mentioned, the ORF-specific 5' primers used for genomic amplification had a first common sequence added thereto, and the ORF-specific 3' primers used for genomic amplification had a second, different, common sequence added thereto, a single set of 5' and 3' primers can be used to amplify all of the probes from the amplifiable ordered set.

Such collections of genome-derived single exon probes can usefully include a plurality of probes chosen for the common attribute of expression in the human heart.

In such defined subsets, typically at least 50, 60, 75, 80, 85, 90 or 95% or more of the probes will be chosen by their expression in the defined tissue or cell type.

The single exon probes of the present invention, as well as fragments of the single exon probes comprising selectively hybridizable portions of the probe ORF, can be used to obtain the full length cDNA that includes the ORF by (i) screening of cDNA libraries; (ii) rapid amplification of cDNA ends ("RACE"); or (iii) other conventional means, as are described, *inter alia*, in Ausubel et al. and Maniatis et al.

It is another aspect of the present invention to provide genome-derived single exon nucleic acid microarrays useful for gene expression analysis, where the term "microarray" has the meaning given in the definitional section of this description, *supra*.

The invention particularly provides genome-derived single-exon nucleic acid microarrays comprising a plurality of probes known to be expressed in human heart. In preferred embodiments, the present invention provides human genome-derived single exon microarrays comprising a plurality of probes drawn from the group consisting of SEQ ID NOS.: 1 - 9,980.

When used for gene expression analysis, the

genome-derived single exon microarrays provide greater physical informational density than do the genome-derived single exon microarrays that have lower percentages of probes known to be expressed commonly in the tested tissue.

5 At a fixed probe density, for example, a given microarray surface area of the defined subset genome-derived single exon microarray can yield a greater number of expression measurements. Alternatively, at a given probe density, the same number of expression measurements can be obtained from

10 a smaller substrate surface area. Alternatively, at a fixed probe density and fixed surface area, probes can be provided redundantly, providing greater reliability in signal measurement for any given probe. Furthermore, with a higher percentage of probes known to be expressed in the

15 assayed tissue, the dynamic range of the detection means can be adjusted to reveal finer levels discrimination among the levels of expression.

Although particularly described with respect to their utility as probes of gene expression, particularly as

20 probes to be included on a genome-derived single exon microarray, each of the nucleic acids having SEQ ID NOS.: 1 - 9,980 contains an open-reading frame, set forth respectively in SEQ ID NOS.: 9,981 - 19,771, that encodes a protein domain. Thus, each of SEQ ID NOS. 1 - 9,980 can be

25 used, or that portion thereof in SEQ ID NOS. 9,981 - 19,771 used, to express a protein domain by standard *in vitro* recombinant techniques. See Ausubel et al. and Maniatis et al.

Additionally, kits are available commercially

30 that readily permit such nucleic acids to be expressed as protein in bacterial cells, insect cells, or mammalian cells, as desired (e.g., HAT™ Protein Expression & Purification System, ClonTech Laboratories, Palo Alto, CA; Adeno-X™ Expression System, ClonTech Laboratories, Palo

35 Alto, CA; Protein Fusion & Purification (pMAL™) System, New

England Biolabs, Beverley, MA)

Furthermore, shorter peptides can be chemically synthesized using commercial peptide synthesizing equipment and well known techniques. Procedures are described, *inter alia*, in Chan et al. (eds.), Fmoc Solid Phase Peptide Synthesis: A Practical Approach (Practical Approach Series, (Paper)), Oxford Univ. Press (March 2000) (ISBN: 0199637245); Jones, Amino Acid and Peptide Synthesis (Oxford Chemistry Primers, No 7) , Oxford Univ. Press (August 1992) (ISBN: 0198556683); and Bodanszky, Principles of Peptide Synthesis (Springer Laboratory), Springer Verlag (December 1993) (ISBN: 0387564314).

It is, therefore, another aspect of the invention to provide peptides comprising an amino acid sequence translated from SEQ ID NOS.: 9,981 - 19,771. Such amino acid sequences are set out in SEQ ID NOS: 19,772 - 29,119. Any such recombinantly-expressed or synthesized peptide of at least 8, and preferably at least about 15, amino acids, can be conjugated to a carrier protein and used to generate antibody that recognizes the peptide. Thus, it is a further aspect of the invention to provide peptides that have at least 8, preferably at least 15, consecutive amino acids.

The following examples are offered by way of illustration and not by way of limitation.

EXAMPLE 1

Preparation of Single Exon Microarrays from ORFs Predicted in Human Genomic Sequence

Bioinformatics Results

All human BAC sequences in fewer than 10 pieces that had been accessioned in a five month period immediately preceding this study were downloaded from

GenBank. This corresponds to ~2200 clones, totaling ~350 MB of sequence, or approximately 10% of the human genome.

After masking repetitive elements using the program CROSS_MATCH, the sequence was analyzed for open
5 reading frames using three separate gene finding programs. The three programs predict genes using independent algorithmic methods developed on independent training sets: GRAIL uses a neural network, GENEFINDER uses a hidden Markoff model, and DICTION, a program proprietary to
10 Genetics Institute, operates according to a different heuristic. The results of all three programs were used to create a prediction matrix across the segment of genomic DNA.

The three gene finding programs yielded a range
15 of results. GRAIL identified the greatest percentage of genomic sequence as putative coding region, 2% of the data analyzed. GENEFINDER was second, calling 1%, and DICTION yielded the least putative coding region, with 0.8% of genomic sequence called as coding region.

20 The consensus data were as follows. GRAIL and GENEFINDER agreed on 0.7% of genomic sequence, GRAIL and DICTION agreed on 0.5% of genomic sequence, and the three programs together agreed on 0.25% of the data analyzed. That is, 0.25% of the genomic sequence was identified by
25 all three of the programs as containing putative coding region.

ORFs predicted by any two of the three programs ("consensus ORFs") were assorted into "gene bins" using two criteria: (1) any 7 consecutive exons within a 25 kb window
30 were placed together in a bin as likely contributing to a single gene, and (2) all ORFs within a 25 kb window were placed together in a bin as likely contributing to a single gene if fewer than 7 exons were found within the 25 kb window.

35

PCR

The largest ORF from each gene bin that did not span repetitive sequence was then chosen for amplification, as were all consensus ORFs longer than 500 bp. This method approximated one exon per gene; however, a number of genes were found to be represented by multiple elements.

Previously, we had determined that DNA fragments fewer than 250 bp in length do not bind well to the amino-modified glass surface of the slides used as support substrate for construction of microarrays; therefore, amplicons were designed in the present experiments to approximate 500 bp in length.

Accordingly, after selecting the largest ORF per gene bin, a 500 bp fragment of sequence centered on the ORF was passed to the primer picking software, PRIMER3 (available online for use at <http://www-genome.wi.mit.edu/cgi-bin/primer/>). A first additional sequence was commonly added to each ORF-unique 5' primer, and a second, different, additional sequence was commonly added to each ORF-unique 3' primer, to permit subsequent reamplification of the amplicon using a single set of "universal" 5' and 3' primers, thus immortalizing the amplicon. The addition of universal priming sequences also facilitates sequence verification, and can be used to add a cloning site should some ORFs be found to warrant further study.

The ORFs were then PCR amplified from genomic DNA, verified on agarose gels, and sequenced using the universal primers to validate the identity of the amplicon to be spotted in the microarray.

Primers were supplied by Operon Technologies (Alameda, CA). PCR amplification was performed by standard techniques using human genomic DNA (Clontech, Palo Alto, CA) as template. Each PCR product was verified by SYBR® green (Molecular Probes, Inc., Eugene, OR) staining of

agarose gels, with subsequent imaging by Fluorimager (Molecular Dynamics, Inc., Sunnyvale, CA). PCR amplification was classified as successful if a single band appeared.

5 The success rate for amplifying ORFs of interest directly from genomic DNA using PCR was approximately 75%. FIG. 5 graphs the distribution of predicted ORF (exon) length and distribution of amplified PCR products, with ORF length shown in red and PCR product length shown in blue
10 (which may appear black in the figure). Although the range of ORF sizes is readily seen to extend to beyond 900 bp, the mean predicted exon size was only 229 bp, with a median size of 150 bp (n=9498). With an average amplicon size of 475 \pm 25 bp, approximately 50% of the average PCR
15 amplification product contained predicted coding region, with the remaining 50% of the amplicon containing either intron, intergenic sequence, or both.

 Using a strategy predicated on amplifying about 500 bp, it was found that long exons had a higher PCR
20 failure rate. To address this, the bioinformatics process was adjusted to amplify 1000, 1500 or 2000 bp fragments from exons larger than 500 bp. This improved the rate of successful amplification of exons exceeding 500 bp, constituting about 9.2% of the exons predicted by the gene
25 finding algorithms.

 Approximately 75% of the probes disposed on the array (90% of those that successfully PCR amplified) were sequence-verified by sequencing in both the forward and reverse direction using MegaBACE sequencer (Molecular
30 Dynamics, Inc., Sunnyvale, CA), universal primers, and standard protocols.

 Some genomic clones (BACs) yielded very poor PCR and sequencing results. The reasons for this are unclear, but may be related to the quality of early draft sequence
35 or the inclusion of vector and host contamination in some

submitted sequence data.

Although the intronic and intergenic material flanking coding regions could theoretically interfere with hybridization during microarray experiments, subsequent
5 empirical results demonstrated that differential expression ratios were not significantly affected by the presence of noncoding sequence. The variation in exon size was similarly found not to affect differential expression ratios significantly; however, variation in exon size was
10 observed to affect the absolute signal intensity (data not shown).

The 350 MB of genomic DNA was, by the above-described process, reduced to 9750 discrete probes, which were spotted in duplicate onto glass slides using
15 commercially available instrumentation (MicroArray GenII Spotter and/or MicroArray GenIII Spotter, Molecular Dynamics, Inc., Sunnyvale, CA). Each slide additionally included either 16 or 32 *E. coli* genes, the average hybridization signal of which was used as a measure of
20 background biological noise.

Each of the probe sequences was BLASTed against the human EST data set, the NR data set, and SwissProt GenBank (May 7, 1999 release 2.0.9).

One third of the probe sequences (as amplified)
25 produced an exact match (BLAST Expect ("E") values less than 1 e^{-100}) to either an EST (20% of sequences) or a known mRNA (13% of sequences). A further 22% of the probe sequences showed some homology to a known EST or mRNA (BLAST E values from 1 e^{-5} to 1 e^{-99}). The remaining 45% of
30 the probe sequences showed no significant sequence homology to any expressed, or potentially expressed, sequences present in public databases.

All of the probe sequences (as amplified) were then analyzed for protein similarities with the SwissProt
35 database using BLASTX, Gish et al., *Nature Genet.* 3:266

(1993). The predicted functional breakdowns of the 2/3 of probes identical or homologous to known sequences are presented in Table 1.

5 Table 1

Function of Predicted ORFs As Deduced From Comparative Sequence Analysis			
Total	V6 chip	V7 chip	Function Predicted from Comparative Sequence Analysis
211	96	115	Receptor
120	43	77	Zinc Finger
30	11	19	Homeobox
25	9	16	Transcription Factor
17	11	7	Transcription
118	57	61	Structural
95	39	56	Kinase
36	18	18	Phosphatase
83	31	52	Ribosomal
45	19	26	Transport
21	17	14	Growth Factor
17	12	5	Cytochrome
50	33	17	Channel

As can be seen, the two most common types of genes were transcription factors and receptors, making up 2.2% and 1.8% of the arrayed elements, respectively.

10

EXAMPLE 2

Gene Expression Measurements From Genome-Derived Single Exon Microarrays

15

The two genome-derived single exon microarrays prepared according to Example 1 were hybridized in a series of simultaneous two-color fluorescence experiments to (1)
5 Cy3-labeled cDNA synthesized from message drawn individually from each of brain, heart, liver, fetal liver, placenta, lung, bone marrow, HeLa, BT 474, or HBL 100 cells, and (2) Cy5-labeled cDNA prepared from message pooled from all ten tissues and cell types, as a control in
10 each of the measurements. Hybridization and scanning were carried out using standard protocols and Molecular Dynamics equipment.

Briefly, mRNA samples were bought from commercial sources (Clontech, Palo Alto, CA and Amersham Pharmacia
15 Biotech (APB)). Cy3-dCTP and Cy5-dCTP (both from APB) were incorporated during separate reverse transcriptions of 1 µg of polyA⁺ mRNA performed using 1 µg oligo(dT)12-18 primer and 2 µg random 9mer primers as follows. After heating to 70°C, the RNA:primer mixture was snap cooled on ice. After
20 snap cooling on ice, added to the RNA to the stated final concentration was: 1X Superscript II buffer, 0.01 M DTT, 100µM dATP, 100 µM dGTP, 100 µM dTTP, 50 µM dCTP, 50 µM Cy3-dCTP or Cy5-dCTP 50 µM, and 200 U Superscript II enzyme. The reaction was incubated for 2 hours at 42°C.
25 After 2 hours, the first strand cDNA was isolated by adding 1 U Ribonuclease H, and incubating for 30 minutes at 37°C. The reaction was then purified using a Qiagen PCR cleanup column, increasing the number of ethanol washes to 5. Probe was eluted using 10 mM Tris pH 8.5.

30 Using a spectrophotometer, probes were measured for dye incorporation. Volumes of both Cy3 and Cy5 cDNA corresponding to 50 pmoles of each dye were then dried in a Speedvac, resuspended in 30 µl hybridization solution containing 50% formamide, 5X SSC, 0.2 µg/µl poly(dA), 0.2
35 µg/µl human c₀t1 DNA, and 0.5 % SDS.

Hybridizations were carried out under a coverslip, with the array placed in a humid oven at 42°C overnight. Before scanning, slides were washed in 1X SSC, 0.2% SDS at 55°C for 5 minutes, followed by 0.1X SSC, 0.2% SDS, at 55°C for 20 minutes. Slides were briefly dipped in water and dried thoroughly under a gentle stream of nitrogen.

Slides were scanned using a Molecular Dynamics Gen3 scanner, as described. Schena (ed.), Microarray Biochip: Tools and Technology, Eaton Publishing Company/BioTechniques Books Division (2000) (ISBN: 1881299376).

Although the use of pooled cDNA as a reference permitted the survey of a large number of tissues, it attenuates the measurement of relative gene expression, since every highly expressed gene in the tissue/cell type-specific fluorescence channel will be present to a level of at least 10% in the control channel. Because of this fact, both signal and expression ratios (the latter hereinafter, "expression" or "relative expression") for each probe were normalized using the average ratio or average signal, respectively, as measured across the whole slide.

Data were accepted for further analysis only when signal was at least three times greater than biological noise, the latter defined by the average signal produced by the *E. coli* control genes.

The relative expression signal for these probes was then plotted as function of tissue or cell type, and is presented in FIG. 6.

FIG. 6 shows the distribution of expression across a panel of ten tissues. The graph shows the number of sequence-verified products that were either not expressed ("0"), expressed in one or more but not all tested tissues ("1" - "9"), and expressed in all tissues tested ("10").

Of 9999 arrayed elements on the two microarrays (including positive and negative controls and "failed" products), 2353 (51%) were expressed in at least one tissue or cell type. Of the gene elements showing significant
5 signal — where expression was scored as "significant" if the normalized Cy3 signal was greater than 1, representing signal 5-fold over biological noise (0.2) — 39% (991) were expressed in all 10 tissues. The next most common class (15%) consisted of gene elements expressed in only a single
10 tissue.

The genes expressed in a single tissue were further analyzed, and the results of the analyses are compiled in FIG. 7.

FIG. 7A is a matrix presenting the expression of
15 all verified sequences that showed expression greater than 3 in at least one tissue. Each clone is represented by a column in the matrix. Each of the 10 tissues assayed is represented by a separate row in the matrix, and relative expression of a clone in that tissue is indicated at the
20 respective node by intensity of green shading, with the intensity legend shown in panel B. The top row of the matrix ("EST Hit") contains "bioinformatic" rather than "physical" expression data — that is, presents the results returned by query of EST, NR and SwissProt databases using
25 the probe sequence. The legend for "bioinformatic expression" (*i.e.*, degree of homology returned) is presented in panel C. Briefly, white is known, black is novel, with gray depicting nonidentical with significant homology (white: E values < 1e-100; gray: E values from 1e-05 to 1e-99; black: E values > 1e-05).
30

As FIG. 7 readily shows, heart and brain were demonstrated to have the greatest numbers of genes that were shown to be uniquely expressed in the respective tissue. In brain, 200 uniquely expressed genes were
35 identified; in heart, 150. The remaining tissues gave the

following figures for uniquely expressed genes: liver, 100; lung, 70; fetal liver, 150; bone marrow, 75; placenta, 100; HeLa, 50; HBL, 100; and BT474, 50.

It was further observed that there were many more
5 "novel" genes among those that were up-regulated in only one tissue, as compared with those that were down-regulated in only one tissue. In fact, it was found that ORFs whose expression was measurable in only a single of the tested tissues were represented in sequencing databases at a rate
10 of only 11%, whereas 36% of the ORFs whose expression was measurable in 9 of the tissues were present in public databases. As for those ORFs expressed in all ten tissues, fully 45% were present in existing expressed sequence databases. These results are not unexpected, since genes
15 expressed in a greater number of tissues have a higher likelihood of being, and thus of having been, discovered by EST approaches.

Comparison of Signal from Known and Unknown Genes

20 The normalized signal of the genes found to have high homology to genes present in the GenBank human EST database were compared to the normalized signal of those genes not found in the GenBank human EST database. The data are shown in FIG. 8.

25 FIG. 8 shows the normalized Cy3 signal intensity for all sequence-verified products with a BLAST Expect ("E") value of greater than $1e-30$ (designated "unknown") upon query of existing EST, NR and SwissProt databases, and shows in blue the normalized Cy3 signal intensity for all
30 sequence-verified products with a BLAST Expect value of less than $1e-30$ ("known"). Note that biological background noise has an averaged normalized Cy3 signal intensity of 0.2.

As expected, the most highly expressed of the
35 ORFs were "known" genes. This is not surprising, since

very high signal intensity correlates with very commonly-expressed genes, which have a higher likelihood of being found by EST sequence.

However, a significant point is that a large
5 number of even the high expressers were "unknown". Since the genomic approach used to identify genes and to confirm their expression does not bias exons toward either the 3' or 5' end of a gene, many of these high expression genes will not have been detected in an end-sequenced cDNA
10 library.

The significant point is that presence of the gene in an EST database is *not* a prerequisite for incorporation into a genome-derived microarray, and further, that arraying such "unknown" exons can help to
15 assign function to as-yet undiscovered genes.

Verification of Gene Expression

To ascertain the validity of the approach described above to identify genes from raw genomic
20 sequence, expression of two of the probes was assayed using reverse transcriptase polymerase chain reaction (RT PCR) and northern blot analysis.

Two microarray probes were selected on the basis of exon size, prior sequencing success, and tissue-specific
25 gene expression patterns as measured by the microarray experiments. The primers originally used to amplify the two respective ORFs from genomic DNA were used in RT PCR against a panel of tissue-specific cDNAs (Rapid-Scan gene expression panel 24 human cDNAs) (OriGene Technologies,
30 Inc., Rockville, MD).

Sequence AL079300_1 was shown by microarray hybridization to be present in cardiac tissue, and sequence AL031734_1 was shown by microarray experiment to be present in placental tissue (data not shown). RT-PCR on these two
35 sequences confirmed the tissue-specific gene expression as

measured by microarrays, as ascertained by the presence of a correctly sized PCR product from the respective tissue type cDNAs.

Clearly, all microarray results cannot, and indeed should not, be confirmed by independent assay methods, or the high throughput, highly parallel advantages of microarray hybridization assays will be lost. However, in addition to the two RT-PCR results presented above, the observation that 1/3 of the arrayed genes exist in expression databases provides powerful confirmation of the power of our methodology – which combines bioinformatic prediction with expression confirmation using genome-derived single exon microarrays – to identify novel genes from raw genomic data.

To verify that the approach further provides correct characterization of the expression patterns of the identified genes, a detailed analysis was performed of the microarrayed sequences that showed high signal in brain.

For this latter analysis, sequences that showed high (normalized) signal in brain, but which showed very low (normalized) signal (less than 0.5, determined to be biological noise) in all other tissues, were further studied. There were 82 sequences that fit these criteria, approximately 2% of the arrayed elements. The 10 sequences showing the highest signal in brain in microarray hybridizations are detailed in Table 2, along with assigned function, if known or reasonably predicted.

Table 2

Function of the Most Highly Expressed Genes Expressed Only in Brain

Microarray Sequence Name	Normal ized Signal	Expressi on Ratio	Homology to EST present in GenBank	Gene Function as described by GenBank
AP000217-1	5.2	+7.7	High	S-100 protein, b-chain, Ca ²⁺ binding protein expressed in central nervous system
AP000047-1	2.3		High	Unknown Function
AC006548-9	1.7		High	Similar to mouse membrane glyco-protein M6, expressed in central nervous system
AC007245-5	1.5		High	Similar to amphiphysin, a synaptic vesicle- associated protein. Ref 21
L44140-4	1.2	+2.0	High	Endothelial actin-binding protein found in nonmuscle filamin

AC004689-9	1.2	+3.5	High	Protein Phosphatase PP2A, neuronal/ downregulates activated protein kinases
AL031657-1	1.2	+3.0	High	Unknown function/ Contains the anhyrin motif, a common protein sequence motif
AC009266-2	1.1	+3.7	Low	Low homology to the Synaptotagmin I protein in rat/present at low levels throughout rat brain
AP000086-1	1.0	+2.7	Low	Unknown, very poor homology to collagen
AC004689-3	1.0		High	Protein Phosphatase PP2A, neuronal/ downregulates activated protein kinases

Of the ten sequences studied by these latter confirmatory approaches, eight were previously known. Of these eight, six had previously been reported to be
5 important in the central nervous system or brain. The exon

giving the highest signal (AP00217-1) was found to be the gene encoding an S100B Ca^{2+} binding protein, reported in the literature to be highly and uniquely expressed in the central nervous system. Heizmann, *Neurochem. Res.* 9:1097
5 (1997).

A number of the brain-specific probe sequences (including AC006548-9, AC009266-2) did not have homology to any known human cDNAs in GenBank but did show homology to rat and mouse cDNAs. Sequences AC004689-9 and AC004689-3
10 were both found to be phosphatases present in neurons (Millward *et al.*, *Trends Biochem. Sci.* 24(5):186-191 (1999)). Two microarray sequences, AP000047-1 and AP000086-1 have unknown function, with AP000086-1 being absent from GenBank. Functionality can now be narrowed
15 down to a role in the central nervous system for both of these genes, showing the power of designing microarrays in this fashion.

Next, the function of the chip sequences with the highest (normalized) signal intensity in brain, regardless
20 of expression in other tissues, was assessed. In this latter analysis, we found expression of many more common genes, since the sequences were not limited to those expressed only in brain. For example, looking at the 20 highest signal intensity spots in brain, 4 were similar to
25 tubulin (AC00807905; AF146191-2; AC007664-4; AF14191-2), 2 were similar to actin (AL035701-2; AL034402-1), and 6 were found to be homologous to glyceraldehyde-3-phosphate dehydrogenase (GAPDH) (AL035604-1; Z86090-1; AC006064-L, AC006064-K; AC035604-3; AC006064-L). These genes are often
30 used as controls or housekeeping genes in microarray experiments of all types.

Other interesting genes highly expressed in brain were a ferritin heavy chain protein, which is reported in the literature to be found in brain and liver (Joshi *et al.*, *J. Neurol. Sci.* 134(Suppl):52-56 (1995)), a result
35

duplicated with the array. Other highly expressed chip
 sequences included a translation elongation factor 1 α
 (AC007564-4), a DEAD-box homolog (AL023804-4), and a Y-
 chromosome RNA-binding motif (Chai *et al.*, *Genomics*
 5 49(2):283-89 (1998)) (AC007320-3). A low homology analog
 (AP00123-1/2) to a gene, DSCR1, thought to be involved in
 trisomy 21 (Down's syndrome), showed high expression in
 both brain and heart, in agreement with the literature
 (Fuentes *et al.*, *Mol. Genet.* 4(10):1935-44 (1995)).

10 As a further validation of the approach, we
 selected the BAC AC006064 to be included on the array.
 This BAC was known to contain the GAPDH gene, and thus
 could be used as a control for the ORF selection process.
 The gene finding and exon selection algorithms resulted in
 15 choosing 25 exons from BAC AC006064 for spotting onto the
 array, of which four were drawn from the GAPDH gene. Table
 3 shows the comparison of the average expression ratio for
 the 4 exons from BAC006064 compared with the average
 expression ratio for 5 different dilutions of a
 20 commercially available GAPDH cDNA (Clontech).

Table 3

Comparison of Expression Ratio, for each tissue, of GAPDH		
	AC006064 (n = 4)	Control (n = 5)
Bone Marrow	-1.81 \pm 0.11	-1.85 \pm 0.08
Brain	-1.41 \pm 0.11	-1.17 \pm 0.05
BT474	1.85 \pm 0.09	1.66 \pm 0.12
Fetal Liver	-1.62 \pm 0.07	-1.41 \pm 0.05
HBL100	1.32 \pm 0.05	2.64 \pm 0.12
Heart	1.16 \pm 0.09	1.56 \pm 0.10
HeLa	1.11 \pm 0.06	1.30 \pm 0.15
Liver	-1.62 \pm 0.22	-2.07 \pm

Lung	-4.95 ± 0.93	-3.75 ± 0.21
Placenta	-3.56 ± 0.25	-3.52 ± 0.43

Each tissue shows excellent agreement between the experimentally chosen exons and the control, again
5 demonstrating the validity of the present exon mining approach. In addition, the data also show the variability of expression of GAPDH within tissues, calling into question its classification as a housekeeping gene and utility as a housekeeping control in microarray
10 experiments.

EXAMPLE 3

Representation of Sequence and Expression Data as a "Mondrian"

15

For each genomic clone processed for microarray as above-described, a plethora of information was accumulated, including full clone sequence, probe sequence within the clone, results of each of the three gene finding
20 programs, EST information associated with the probe sequences, and microarray signal and expression for multiple tissues, challenging our ability to display the information.

Accordingly, we devised a new tool for visual
25 display of the sequence with its attendant annotation which, in deference to its visual similarity to the paintings of Piet Mondrian, is hereinafter termed a "Mondrian". FIGS. 3 and 4 present the key to the information presented on a Mondrian.

30 FIG. 9 presents a Mondrian of BAC AC008172 (bases 25,000 to 130,000 shown), containing the carbamyl phosphate synthetase gene (AF154830.1). Purple background within the region shown as field 81 in FIG. 3 indicates all 37 known

exons for this gene.

As can be seen, GRAIL II successfully identified 27 of the known exons (73%), GENEFINDER successfully identified 37 of the known exons (100%), while DICTION
5 identified 7 of the known exons (19%).

Seven of the predicted exons were selected for physical assay, of which 5 successfully amplified by PCR and were sequenced. These five exons were all found to be from the same gene, the carbamyl phosphate synthetase gene
10 (AF154830.1).

The five exons were arrayed, and gene expression measured across 10 tissues. As is readily seen in the Mondrian, the five chip sequences on the array show identical expression patterns, elegantly demonstrating the
15 reproducibility of the system.

FIG. 10 is a Mondrian of BAC AL049839. We selected 12 exons from this BAC, of which 10 successfully sequenced, which were found to form between 5 and 6 genes. Interestingly, 4 of the genes on this BAC are protease
20 inhibitors. Again, these data elegantly show that exons selected from the same gene show the same expression patterns, depicted below the red line. From this figure, it is clear that our ability to find known genes is very good. A novel gene is also found from 86.6 kb to 88.6 kb,
25 upon which all the exon finding programs agree. We are confident we have two exons from a single gene since they show the same expression patterns and the exons are proximal to each other. Backgrounds in the following colors indicate a known gene (top to bottom):
30 red = kallistatin protease inhibitor (P29622);
purple = plasma serine protease inhibitor (P05154);
turquoise = α 1 anti-chymotrypsin (P01011); mauve = 40S ribosomal protein (P08865). Note that chip sequence 8 and 12 did not sequence verify.

35

EXAMPLE 4Genome-Derived Single Exon Probes Useful For Measuring
Human Gene Expression

5

The protocols set forth in Examples 1 and 2, *supra*, were applied to additional human genomic sequence as it became newly available in GenBank to identify unique exons in the human genome that could be shown to be
10 expressed at significant levels in heart tissue.

These unique exons are within longer probe sequences. Each probe was completely sequenced on both strands prior to its use on a genome-derived single exon microarray; sequencing confirms the exact chemical
15 structure of each probe. An added benefit of sequencing is that it placed us in possession of a set of single base-incremented fragments of the sequenced nucleic acid, starting from the sequencing primer 3' OH. (Since the single exon probes were first obtained by PCR amplification
20 from genomic DNA, we were of course additionally in possession of an even larger set of single base incremented fragments of each of the 9,980 single exon probes, each fragment corresponding to an extension product from one of the two amplification primers.)

25 The structures of the 9,980 unique single exon probes are clearly presented in the Sequence Listing as SEQ ID Nos.: 1 - 9,980. The 16 nt 5' primer sequence and 16 nt 3' primer sequence present on the amplicon are not included in the sequence listing. The sequences of the exons
30 present within each of these probes is presented in the Sequence Listing as SEQ ID Nos.: 9,981 - 19,771, respectively. It will be noted that some amplicons have more than one exon, some exons are contained in more than one amplicon.

35 As detailed in Example 2, expression was

demonstrated by disposing the amplicons as single exon probes on nucleic acid microarrays and then performing two-color fluorescent hybridization analysis; significant expression is based on a statistical confidence that the
5 signal is significantly greater than negative biological control spots. The negative biological control is formed from spotted DNA sequences from a different species. Here, 32 sequences from E.Coli were spotted in duplicate to give a total of 64 spots.

10 For each hybridisation (each slide, each colour) the median value of the signal from all of the spots is determined. The normalised signal value is the arithmetic mean of the signal from duplicate spots divided by the population median.

15 Control spots are eliminated if there is more than a five-fold difference between each one of the duplicate spots raw signals.

The median of the signal from the remaining control spots is calculated and all subsequent calculations
20 are done with normalised signals.

Control spots having a signal of greater than median + 2.4 (the value 2.4 is roughly 12 times the observed standard deviation of control spot populations) are eliminated. Spots with such high signals are considered
25 to be "outliers".

The mean and standard deviation of the modified control spot populations are calculated.

The mean + 3x the standard deviation (mean + (3*SD)) is used as the signal threshold qualifier for that
30 particular hybridisation. Thus, individual thresholds are determined for each channel and each hybridisation.

This means that, assuming that the data is distributed normally, there is a 99% confidence that any signal exceeding the threshold is significant.

35 The probes and their expression data are

presented in Table 4, set forth respectively in Example 5. Example 5 presents the subset of probes that is significantly expressed in the human heart and thus presents the subset of probes that was recognized to be
5 useful for measuring expression of their cognate genes in human heart tissue.

The sequence of each of the exon probes identified by SEQ ID NOS.: 9,981 - 19,771 was individually used as a BLAST (or, for SWISSPROT, BLASTX) query to
10 identify the most similar sequence in each of dbEST, SwissProt (BLASTX), and NR divisions of GenBank. Because the query sequences are themselves derived from genomic sequence in GenBank, only nongenic hits from NR were scored.

15 The smallest in value of the BLAST (or BLASTX) expect ("E") scores for each query sequence across the three database divisions was used as a measure of the "expression novelty" of the probe's ORF. Table 4 is sorted in descending order based on this measure, reported as
20 "Most Similar (top) Hit BLAST E Value". Those sequences for which no "Hit E Value" is listed are those exons which were found to have no similar sequences.

As sorted, Table 4 thus lists its respective probes (by "AMPLICON SEQ ID NO.:" and additionally by the
25 SEQ ID NO.: of the exon contained within the probe:"EXON SEQ ID NO.:") from least similar to sequences known to be expressed (i.e., highest BLAST E value), at the beginning of the table, to most similar to sequences known to be expressed (i.e., lowest BLAST E value), at the bottom of
30 the table.

Table 4 further provides, for each listed probe, the accession number of the database sequence that yielded the "Most Similar (top) Hit BLAST E Value", along with the name of the database in which the database sequence is
35 found ("Top Hit Database Source").

Table 4 further provides SEQ ID NOS.

corresponding to the predicted amino acid sequences where they have been determined for the probe and exon nucleotide sequences. These are set out as PEPTIDE SEQ ID NOS.:. The
5 peptide sequences for a given exon are predicted as follows: Since each chip exon is a consensus sequence drawn from predictions from various exon finding programs (i.e. Grail, GeneFinder and GenScan), the multiple initial ORFs are first determined in a uniform way according to each
10 prediction. In particular, the reading frame for predicting the first amino acid in the peptide sequence always starts with the first base of any codon and ends with the last base of non-termination codon. Next, for each strand of the exon, initial ORFs are merged into one or more final ORFs
15 in an exhaustive process based on the following criteria: 1) the merging ORFs must be overlapping, and 2) the merging ORFs must be in the same frame.

The Sequence Listing, which is a superset of all of the data presented in Table 4, further includes, for
20 each probe, the most similar hit, with accession number and BLAST E value, from the each of the three queried databases.

Table 4 further lists, for each probe, a portion of the descriptor for the top hit ("Top Hit Descriptor") as
25 provided in the sequence database. For those ORFs that are similar in sequence, but nonidentical to known sequences (e.g., those with BLAST E values between about $1e-05$ and $1e-100$), the descriptor reveals the likely function of the protein encoded by the probe's ORF.

30 Using BLAST E value cutoffs of $1e-05$ (i.e., 1×10^{-5}) and $1e-100$ (i.e., 1×10^{-100}) as evidence of similarity to sequences known to be expressed is of course arbitrary: in Example 2, *supra*, a BLAST E value of $1e-30$ was used as the boundary when only two classes were to be defined for
35 analysis (unknown, $>1e-30$; known $<1e-30$) (see also FIG. 8).

Furthermore, even when the "Most Similar (Top) Hit BLAST E Value" is low, e.g., less than about $1e-100$ — which is probative evidence that the query sequence has previously been shown to be expressed — the top hit is highly unlikely
5 exactly to match the probe sequence.

First, such expression entries typically will not have the intronic and/or intergenic sequence present within the single exon probes listed in the Table. Second, even the ORF itself is unlikely in such cases to be present
10 identically in the databases, since most of the EST and mRNA clones in existing databases include multiple exons, without any indication of the location of exon boundaries.

As noted, the data presented in Table 4 represent a proper subset of the data present within the attached
15 sequence listing. For each amplicon probe (SEQ ID NOs.: 1 - 9,980) and probe exon (SEQ ID NOs.: 9,981 - 19,771, respectively), the sequence listing further provides, through iterated annotation fields <220> and <223>:

. . (a) the accession number of the BAC from which
20 the sequence was derived ("MAP TO"), thus providing a link to the chromosomal map location and other information about the genomic milieu of the probe sequence;

(b) the most similar sequence provided by BLAST query of the EST database, with accession number and BLAST
25 E value for the "hit";

(c) the most similar sequence provided by BLAST query of the GenBank NR database, with accession number and BLAST E value for the "hit"; and

(d) the most similar sequence provided by BLASTX
30 query of the SWISSPROT database, with accession number and BLAST E value for the "hit".

EXAMPLE 5

35 Genome-Derived Single Exon Probes Useful For Measuring